

# SCIENTIFIC REPORTS



OPEN

## The evolutionary life cycle of the polysaccharide biosynthetic gene cluster based on the *Sphingomonadaceae*

Received: 10 November 2016

Accepted: 21 March 2017

Published: 21 April 2017

Mengmeng Wu<sup>1,\*</sup>, Haidong Huang<sup>2,\*</sup>, Guoqiang Li<sup>1</sup>, Yi Ren<sup>3</sup>, Zhong Shi<sup>1</sup>, Xiaoyan Li<sup>2</sup>, Xiaohui Dai<sup>1</sup>, Ge Gao<sup>1</sup>, Mengnan Ren<sup>2</sup> & Ting Ma<sup>1</sup>

Although clustering of genes from the same metabolic pathway is a widespread phenomenon, the evolution of the polysaccharide biosynthetic gene cluster remains poorly understood. To determine the evolution of this pathway, we identified a scattered production pathway of the polysaccharide sanxan by *Sphingomonas sanxanigenens* NX02, and compared the distribution of genes between sphingans-producing and other *Sphingomonadaceae* strains. This allowed us to determine how the scattered sanxan pathway developed, and how the polysaccharide gene cluster evolved. Our findings suggested that the evolution of microbial polysaccharide biosynthesis gene clusters is a lengthy cyclic process comprising cluster 1 → scatter → cluster 2. The sanxan biosynthetic pathway proved the existence of a dispersive process. We also report the complete genome sequence of NX02, in which we identified many unstable genetic elements and powerful secretion systems. Furthermore, nine enzymes for the formation of activated precursors, four glycosyltransferases, four acyltransferases, and four polymerization and export proteins were identified. These genes were scattered in the NX02 genome, and the positive regulator SpnA of sphingans synthesis could not regulate sanxan production. Finally, we concluded that the evolution of the sanxan pathway was independent. NX02 evolved naturally as a polysaccharide producing strain over a long-time evolution involving gene acquisitions and adaptive mutations.

A gene cluster is a set of functionally related genes located in close physical proximity in a genome, and an operon, a more structured instance of a cluster, refers to a set of genes under common regulatory control<sup>1,2</sup>. Clustering of genes involved in the same metabolic pathway is a widespread phenomenon<sup>3,4</sup>. The genes of certain biosynthetic pathways for microbial polysaccharides<sup>5,6</sup>, such as cellulose<sup>7</sup>, alginate<sup>8</sup>, succinoglycan<sup>9</sup>, sphingans<sup>10,11</sup> and xanthan<sup>12</sup>, form clusters and operons in the genome. Some models suggest a tendency for genes to cluster. For instance, the selfish operons model postulates that genes organized into a cluster can propagate by vertical transmission and horizontal transfer<sup>13–15</sup>, which might represent an instinct for self-preservation. The co-regulation model states that the formation of an operon promotes the production of gene products in equal measures<sup>16</sup>. The alternative explanation for cluster formation is known as the protein immobility model, which suggests that clustered genes produce local clusters of enzymes, and the physical proximity of the enzymes minimize the steady state level of reaction step intermediates and thus conserves energy and material required for growth and maintenance<sup>17</sup>. Thus, in the biosynthesis of a polysaccharide, genes might form clusters and operons. However, the genome structures within operons are unstable and shuffling of a genome structure is virtually neutral over long-term evolution<sup>18</sup>. The changes to operon structures can represent identical structure, similar structures (translocation, deletion, and insertion), destructed structures and complete loss<sup>18</sup>. In addition, Morgan *et al.*<sup>1</sup> determined that the evolutionary life cycle of operons occurs via the mechanisms of operon formation and death by gene insertion, deletion or rearrangement. To date, there has been no report of the evolutionary

<sup>1</sup>Key Laboratory of Molecular Microbiology and Technology, Ministry of Education, College of Life Sciences, Nankai University, Tianjin, China. <sup>2</sup>College of Agronomy & Resources and Environment, Tianjin Agricultural University, Tianjin, China. <sup>3</sup>Shanghai Majorbio Bio-pharm Biotechnology Limited Company, Shanghai, China. \*These authors contributed equally to this work. Correspondence and requests for materials should be addressed to T.M. (email: tingma@nankai.edu.cn)

history of genes or clusters related to polysaccharide biosynthesis. Evolutionary research on these pathways is very difficult because of the clustering of responsible polysaccharide biosynthesis genes and exopolysaccharide (EPS) mutated strains with disrupted clusters during long-term evolution have received little research attention. Furthermore, the absence of recognized scattered polysaccharide synthetic pathway also makes exploring its evolution difficult.

The genus *Sphingomonas* belongs to the family *Sphingomonadaceae*<sup>19</sup>. Subsequently, the genus description was amended<sup>19–22</sup>. Based on phylogenetic analysis, polyamine patterns and fatty acid profiling, the genus *Sphingomonas* was subdivided into four genera: *Sphingomonas sensu stricto* and three new genera, *Sphingobium*, *Novosphingobium* and *Sphingopyxis*. The *Sphingomonadaceae* have great potential for biotechnological applications in bioremediation, the degradation of refractory contaminants and the production of valuable biopolymers. The biopolymers synthesized by strains of the *Sphingomonadaceae* have similar but not identical structures and are termed sphingans, including gellan, welan, rhamsan, S-88, S-7, S-198 and S-657<sup>10,11</sup>. Sphingans possess a similar linear repeating tetrasaccharide of [→4) α-L-Rha (1→3) β-D-Glc (1→4) β-D-GlcA (1→4) β-D-Glc (→1] (Supplementary Fig. S1)<sup>10</sup>. Sphingans are used in food, pharmaceutical, and other industries as a stabilizing, thickening, emulsifying, and gelling agents because of their excellent rheological characteristics<sup>10,11</sup>.

The general scheme of sphingan biosynthesis follows the Wzx/Wzy-dependent pathway as described for *Escherichia coli* group I and IV polysaccharide<sup>23,24</sup>. Biosynthetic pathways for sphingan S-88<sup>25</sup>, S-7<sup>26</sup>, gellan<sup>10,27</sup>, diutan<sup>28</sup> and welan<sup>11,29</sup> have been identified. Genes related to the assembly, polymerization and export of tetrasaccharide repeat units are clustered with same arrangement and orientation, and named the *spn* cluster<sup>11</sup>. Biochemical analysis indicated that the function of glycosyltransferase (GT) *spnB*<sup>26</sup> and *spnK*<sup>30</sup>. Protein *spnC* and *spnE* potentially involved in secretion and chain length determination have also been analyzed biochemically<sup>31</sup>. The detailed enzymological functions of GT SpnL/Q, polymerase SpnG, and flippase SpnS remain unclear. In the sphingan gene cluster, the function of *spnI/J/F/M/N* are unknown, but might be related indirectly to sphingan biosynthesis<sup>10,32</sup>. The putative lyase SpnR might excise the mature polysaccharide chain from the outer membrane. Regulatory protein SpnA with regions homologous to sensor kinases and response regulator proteins<sup>11</sup> is located far away from the sphingan cluster, thus its regulatory mechanism requires further research.

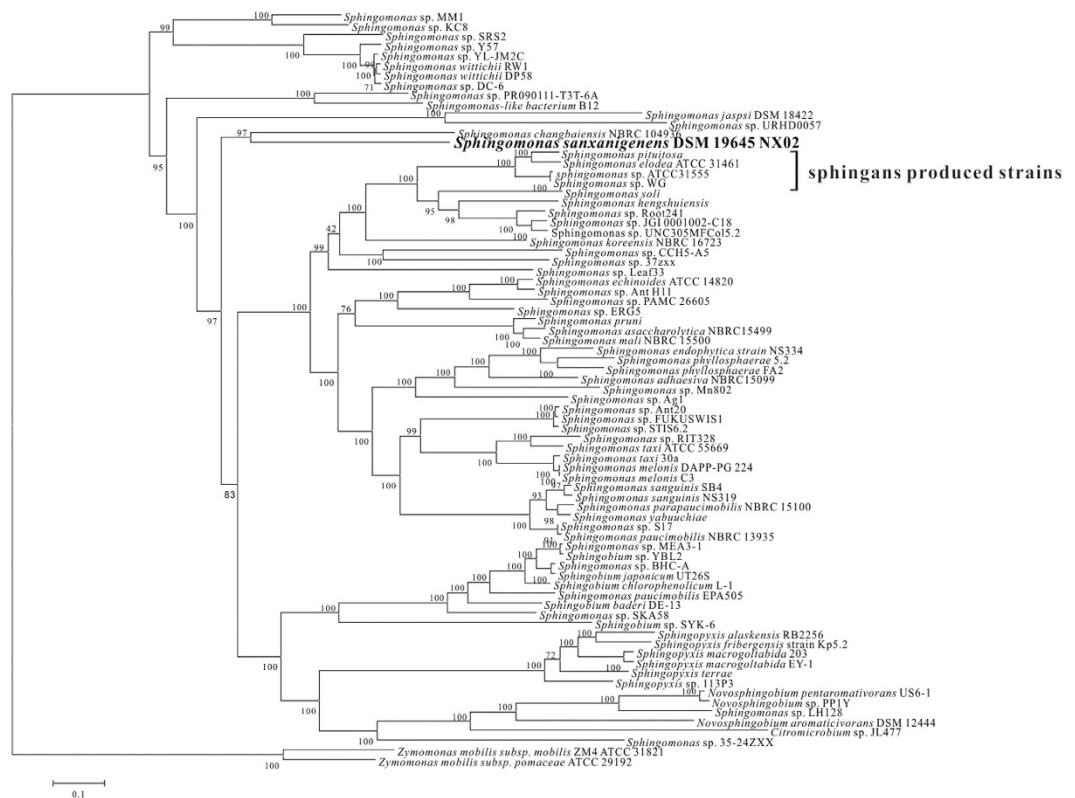
Newly isolated *Sphingomonas sanxanigenens* NX02 (DSM 19645), a Gram-negative and rod-shaped bacterium<sup>33</sup>, synthesizes a novel extracellular polymer sanxan, comprising a tetrasaccharide repeat unit of [→4) β-D-Man (1→4) β-D-GlcA (1→3) α-L-Rha (1→3) β-D-Glc (→1]<sup>34</sup>, which was very different from other sphingans (Supplementary Fig. S1). Sanxan has excellent thickening, shear thinning, gelling, and emulsification properties<sup>35</sup>, and has been used in China for years as drilling mud and as a thickening agent to recover petroleum by water flooding. In this study, we proposed a hypothetical evolutionary model for the distribution of genes for microbial polysaccharide biosynthesis, based on *Sphingomonas sanxanigenens* NX02 and the *Sphingomonadaceae*, comprising a lengthy cyclic process: cluster 1 → scatter → cluster 2. The genes responsible for sanxan biosynthesis prove the existence of a dispersive process. This evolutionary process was unstable and dynamic, and cluster 1 and 2 contained non-identical genes and rearrangement. In addition, we report the complete genome sequence of NX02 and the identification of a scattered pathway for sanxan biosynthesis, which partially explains sanxan's specific structure. This is the first report that a scattered pathway could synthesize a polysaccharide. The analysis of the genome provided insights into its large genome size.

## Results and Discussion

**Genome features of *S. sanxanigenens* NX02.** The complete genome of *Sphingomonas sanxanigenens* NX02 comprises a circular 6,205,897-bp chromosome and a 374,401-bp plasmid, with G + C contents of 66.8% and 64.9%, respectively. The genome of strain NX02 is the largest in the *Sphingomonadaceae* submitted to GenBank. The chromosome is predicted to contain 5,619 protein coding genes (CDSs), with an average size of 1,011 bp, nine rRNA operons and 58 tRNA genes. The plasmid contains 356 predicted CDSs with an average size of 903 bp, giving a coding intensity of 87.9%.

The phylogenetic relationships of *S. sanxanigenens* NX02 with 77 sequenced strains of the *Sphingomonadaceae* (submitted before April 2016) based on all core genes were analysed (Fig. 1). Most of these 77 sequenced strains, such as *Sphingomonas wittichii* RW1<sup>36</sup>, *Sphingobium chlorophenicum* L-1<sup>37</sup> and *Sphingobium japonicum* UT26S<sup>38</sup>, were reported to degrade persistent complex compound<sup>39</sup>. Besides NX02, only four sequenced strains were reported to produce polysaccharide like sphingans: *Sphingomonas elodea* ATCC 31461 (Gellan)<sup>40</sup>, *Sphingomonas* sp. ATCC31555 (Welan)<sup>41</sup>, *Sphingomonas* sp. WG (Welan)<sup>42</sup> and *Sphingomonas puititosa* (PS-EDIV)<sup>43</sup>. In addition, there was no EPS formation reported for the 72 non-sphingan producing strains. NX02 formed a separate evolutionary clade with *Sphingomonas changbaiensis* (isolated from Changbai mountains, China<sup>44</sup>), and the differentiation between them happened at an early stage of evolution. Despite being biopolymer-producing strain, NX02 was evolutionarily distant from the sphingan-producing strains.

**Genomic islands and horizontal gene transfer (HGT).** The NX02 genome has undergone a number of HGT events assisted by phages and transposons. In all, 327 insertion sequences (ISs, E value < 1.00<sup>e-20</sup>), representing 93.8 kb (Supplementary Table S1), and 104 transposases were annotated in the NX02 genome. The most frequently identified were ISMd17 (28 copies; 706 bp, 482 bp, 139 bp, and 212 bp appeared seven times, respectively) which originated from *Methylobacterium dichlormethanicum*. Twelve phages (215.6 kb) were found in the genome, nine of which were on the chromosome and three on the plasmid (Supplementary Table S2). An intact prophage was found in the chromosome (4076698–4105468), comprising 28.7 kb, with 35 CDS. In NX02, 45 gene islands (GIs), comprising 298 kb and 308 CDS, were identified, including transposases, transcriptional regulators and membrane proteins (Supplementary Table S3, Supplementary discussion). The distributions of ISs, prophages and GIs in the genome are shown in Fig. 2. Thus, NX02's genome is very unstable and active, and the many ISs,



**Figure 1.** Neighbour-Joining phylogenetic tree of 77 *Spingomonad* genomes constructed from concatenated nucleotide sequences of universally conserved genes using the Mega 6 tool. The numbers for the interior branches are bootstrap percentages. The scale bar indicates the number of substitutions per site. The four sphingan-producing strains are *S. pituitosa*, *S. elodea* ATCC 31461, *Spingomonas* sp. ATCC31, and *Spingomonas* sp. WG.

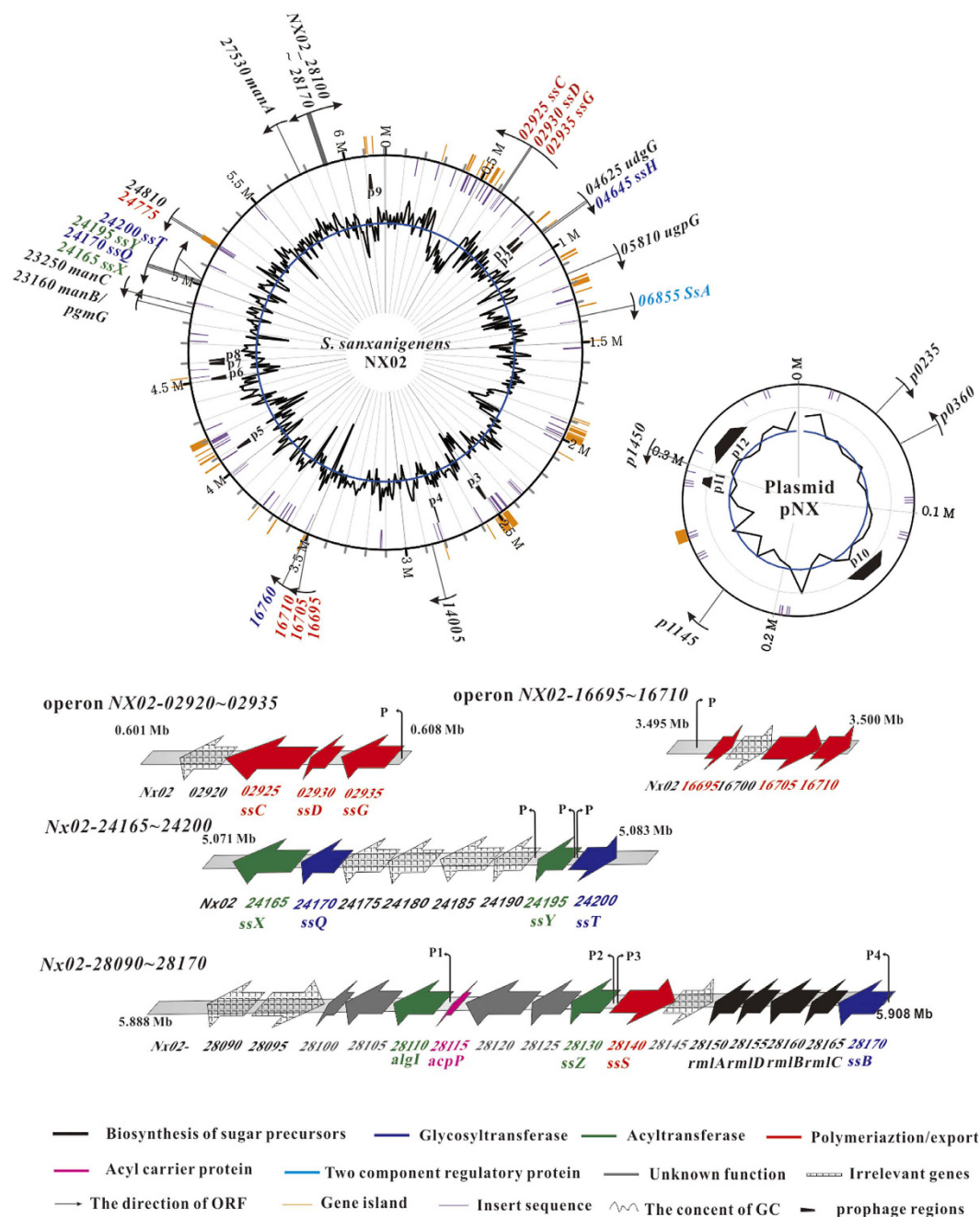
transposases, GIs, and prophages would result in gene duplication, rearrangement, and loss, which accelerate genome evolution.

**Secretion systems and DNA uptake.** Five typical secretion systems (SS) were found in the NX02 genome, including one T1SS, one T2SS, one T3SS, four T4SS, and one T6SS (Supplementary Fig. S2, Supplementary discussion). By comparison, *Spingomonas* sp. ATCC 31555 contains one set of T1SS–T4SS, and *S. elodea* ATCC31461 and *S. wittichii* RW1 only have one T1SS, T2SS and T4SS. The T1SS of NX02 comprised a single copy protein TolC (outer membrane protein, OMP), 11 copies of HlyD as the MFP (membrane fusion protein) and two copies of ABC (ATP-binding cassette, which showed low homology with the HlyB protein, but was near the MFP) and might be responsible for drug resistance. The T3SS comprised 11 gene products from NX02\_15795 to NX02\_15875, which all shared higher homology and were classified as T3SS components form *Spingomonas* sp. SKA58. This cluster was considered a “pathogenicity island”. Four sets of T4SSs are present in NX02, namely T4SS-1 (NX02\_p0495 to NX02\_p0580), T4SS-2 (NX02\_09735 to NX02\_09825), T4SS-3 (NX02\_11790 to NX02\_11845), and T4SS-4 (NX02\_19725 to NX02\_19775). T4SS-1 is on the plasmid, while T4SS-1, -2 and -3 are on the chromosome. The T6SS of NX02 comprises 15 genes in the *imp* operon, including *vgrG*, *hcp*, *vasU*, and *clpV*. Interestingly, 12 of the NX02 T6SS genes are most similar to *imp* genes from *Spingomonas* sp. S17.

Naturally competent bacteria use certain proteins to take up DNA. Parts of this common competence system share homology with proteins that are involved in the assembly of type IV pili and type II secretion systems, and form a structure that spans the cell envelope partially<sup>45,46</sup>. Interestingly, NX02 has a T2SS and four sets of T4SS, suggesting that NX02 has greater capacity to take up exogenous DNA.

Thus, the secretion systems of NX02 probably play a critical role in HGT, permitting adaption to the environment, and driving bacterial. These multiple TSSs might explain the huge genome. In addition, the ISs, transposases, GIs, and prophages could also enlarge the genome by multiple gene duplication.

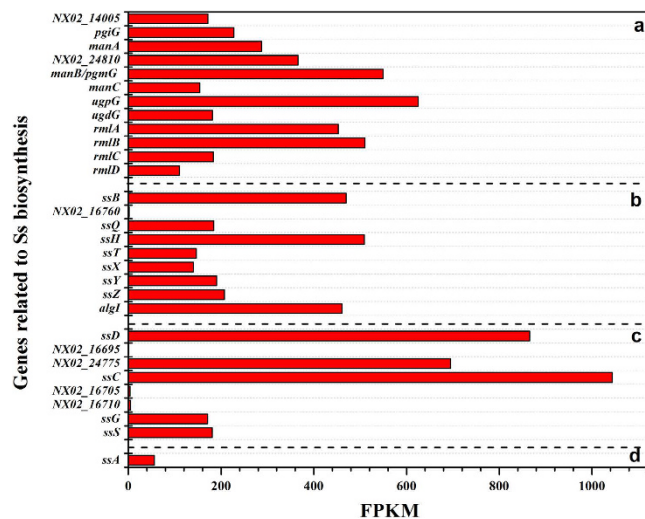
**The identification of the biosynthetic pathway of the biopolymer sanxan.** The biosynthetic pathways of sphingans are similar to those of group 1 and 4 capsule polysaccharides from *E. coli*<sup>23,24</sup>. The sanxan biosynthetic pathway also comprises a multi-step Wzx/Wzy-dependent process<sup>11,23,24</sup>, which is divided into three sequential steps: (a) sugar-activated precursors are synthesized simultaneously; (b) tetrasaccharide repeat units are assembled into the inner membrane; and (c) the repeat units are polymerized and exported through the outer membrane. The detailed process and enzymes involved in sanxan synthesis are as follows.



**Figure 2.** Genome features of *Spingomonas sanxanigenens* NX02. The genome plot showing (from the outside to the centre) genes related to sanxan biosynthesis (different colours or symbols stand for genes with different functions), gene islands, genome size, insert sequence, prophage, and the GC content. The features of the four clusters related to sanxan synthesis are also shown.

**Genes and enzymes involved in the formation of nucleotide sugar precursors.** The tetrasaccharide repeat units of sanxan are synthesized from activated UDP-D-glucose, UDP-D-glucuronic acid, GDP-D-mannose, and dTDP-L-Rhamnose (Supplementary discussion). The *pgmG* gene, encoding a phosphoglucomutase (EC 5.4.2.2), is represented by *NX02\_14005* and *NX02\_23160*<sup>47</sup>. Their FPKM (Fragments Per Kilobase of exon per Million mapped reads) values were both high (Fig. 3). The *NX02\_14005* disruptant had an S<sup>+</sup> phenotype like the wild-type strain, while the *NX02\_23160* mutant showed an S<sup>-</sup> phenotype in NKS medium (figure not shown). The sanxan yield of the *NX02\_23160*-deficient strain in fermentation broth observably decreased, and its overexpression resulted in a 17 ± 0.3% increase sanxan production<sup>47</sup>. Therefore, *NX02\_23160* encodes the main phosphoglucomutase, PgmG, while *NX02\_14005* might be an alternate protein.

The *ugpG* gene, encoding a UDP-glucose pyrophosphorylase (EC 2.7.7.9), which catalyses the reversible conversion of glucose-1-phosphate and UTP into precursor UDP-D-glucose and diphosphate<sup>48</sup>, is *NX02\_05810*. The *ugdG* gene, encoding a UDP-glucose dehydrogenase (EC 1.1.1.22), which converts UDP-glucose into the



**Figure 3.** FPKM (Fragments Per Kilobase of exon per Million mapped reads) values of genes related to sanxan biosynthesis on NK medium. “a” stands for genes related to precursor biosynthesis; “b” stands for genes responsible for the assembly of the repeat units; “c” indicates genes related to the polymerization and export of sanxan polysaccharides; “d” is the *ssA* gene.

activated precursor UDP-D-glucuronic acid, is *NX02\_04625*<sup>49</sup>. The enzymes required for the synthesis of precursor dTDP-L-Rhamnose are TDP-glucose pyrophosphorylase (*RmlA*, EC 2.7.7.24, *NX02\_28150*), dTDP-glucose 4,6-dehydratase (*RmlB*, EC 4.2.1.46, *NX02\_28160*), dTDP-4-dehydrorhamnose 3,5-epimerase (*RmlC*, EC 5.1.3.13, *NX02\_28165*), and dTDP-4-dehydrorhamnose reductase (*RmlD*, EC 1.1.1.133, *NX02\_28155*). The wild-type strain deleted for the *rml* cluster was not obtained despite screening a large number of mutants. However, the mutant strain could be obtained in a glucosyl-isoprenylphosphate transferase-deficient strain.

The activated precursor GDP-D-mannose is synthesized by mannose-6-phosphate isomerase (*ManA*, EC 5.3.1.8, *NX02\_27530*), phosphomannomutase (*ManB*, EC 5.4.2.8, bifunctional gene *NX02\_28160*), and Mannose-1-phosphate guanylyltransferase (*ManC*, EC 2.7.7.22, *NX02\_23250*). These related genes were scattered among the genome. Identity or similarity analysis between proteins related to precursor synthesis from *S. sanxanigenens* and *S. elodea* ATCC 31461, *Sphingomonas* sp. ATCC 31555, and *Sphingomonas* sp. ATCC 53159 are shown in Table 1.

**Genes and enzymes involved in the assembly of the tetrasaccharide repeat unit.** The synthesis of activated precursors was followed by the formation of the tetrasaccharide repeat unit by sequential transfer of the sugar and acyl donors to an activated lipid carrier by glycosyltransferases (GTs) and acyltransferases (ATs), according to the structure of sanxan (Lipid-P-P ←  $\beta$ -D-Glc ←  $\alpha$ -L-Rha ←  $\beta$ -D-GlcA ←  $\beta$ -D-Man, Supplementary Fig. S1). The tetrasaccharide repeat units were assembled on a lipid carrier comprising the C<sub>55</sub>-isoprenylphosphate carrier (IP), which was similar to the group I capsular polysaccharide in *E. coli*<sup>23,24</sup>.

The priming glycosyltransferase (*SsB*) annotated as glucose-1-isoprenylphosphate transferase, which transfers glucose-1-phosphate from UDP-glucose to the lipid carrier IP, was encoded by *NX02\_28170*. The *ssB* gene is located upstream of the *rml* cluster and shares a common promoter (Fig. 2). The *ssB*-deficient strain showed an Ss<sup>-</sup> phenotype and was complemented by plasmid pBBR*ssB* or pBBR*gelB*. The genotype and phenotype of the mutant and complementing strains are shown in Fig. 4. Protein *SsB* has four predicted N-terminal transmembrane regions and one at Leu<sub>280</sub> to Val<sub>301</sub>; its C-terminus is predicted to be cytoplasmic. *SsB* is homologous to *GelB* from *S. elodea* ATCC31461 (41.6%)<sup>28</sup>, *WelB* from *Sphingomonas* sp. ATCC 31555 (43.4%), and *SpsB* from *Sphingomonas* sp. ATCC 53159 (44.1%), Table 1. In addition, another gene, *NX02\_16760*, was also predicted as glucose-1-isoprenylphosphate transferase, however its FPKM value was low, and the mutant strain had the Ss<sup>+</sup> phenotype with unchanged sanxan production (Supplementary Fig. S3). Thus, under most circumstances, *NX02\_16760* is probably irrelevant to sanxan biosynthesis.

In the biosynthesis of gellan, welan, diutan, or S-88, genes related to assembly, polymerization, and export are clustered with almost the same arrangement<sup>11</sup>. However, similar gene clusters were not found within 50 kb of *ssB*. Other GTs were located in separate loci in the genome. According to the CAZY database and gene annotation, 33 GTs were found in the chromosome. The purified yields of strains deficient for these GTs and their FPKM values are shown in the Supplementary discussion and Supplementary Fig. S3. Among of all 33 genes, *NX02\_24170*, *NX02\_24200* and *NX02\_04645* were special (Supplementary discussion).

Deletion of *NX02\_24170* and *NX02\_24200* failed in the wild-type strain, while a mutant could be obtained in strain *NX02* ( $\Delta$ *ssB*). Complementation experiments showed that strain *NX02* ( $\Delta$ *ssB*,  $\Delta$ *24170*) had an Ss<sup>+</sup> phenotype when plasmid pBBR*ssrgelQ* was transferred (Fig. 4). *NX02\_24170* was homologous to *gelQ* (only 22.6% identity, Table 1) and it was named as *ssQ*. *SsQ* is demonstrated to be the second glycosyltransferase that transfers rhamnose from dTDP-L-rhamnose to IPP ← glucose. Strain *NX02* ( $\Delta$ *ssB*,  $\Delta$ *24200*) could be complemented by plasmid pBBR*ssrgelL*. *GelL* catalyses the addition of  $\beta$ -D-glucose to  $\beta$ -D-glucuronic acid in gellan

Sanxan	Amino acids	Predicted function	% Identity (% Similarity)			Accession no.
			Gellan	Welan	Diutan	
<b>Biosynthesis of nucleotide-sugar precursors</b>						
<i>manB/pgmG</i>	460	Phosphoglucomutase/phosphomannomutase	80.3 (87.7)	80.1 (87.2)	—	AHE56248.1
<i>manA</i>	277	Mannose-6-phosphate isomerase	71.7 (80.4)	71.5 (80.4)	—	AHE57091.1
<i>manC</i>	342	Mannose-1-phosphate guanylyltransferase	54.8 (66.7)	63.2 (76.8)	—	AHE56266.1
<i>ugpG</i>	288	Glucose-1-phosphate uridylyltransferase	76.2 (86.1)	74.8 (86.9)	—	AHE52897.1
<i>ugdG</i>	454	UDP-glucose-6-dehydrogenase	70.8 (80.4)	70.5 (80.4)	—	AHE52667.1
<i>rmlA</i>	288	Glucose-1-phosphate thymidyltransferase	60.3 (77.4)	60.6 (76.4)	61.6 (77.7)	AHE57212.1
<i>rmlB</i>	356	dTDP-glucose 4,6-dehydratase	63.1 (74.3)	63.4 (74.6)	63.1 (74.3)	AHE57214.1
<i>rmlC</i>	186	dTDP-4-dehydrorhamnose 3,5-epimerase	48.2 (62.6)	49.2 (62.6)	48.2 (63.1)	AHE57215.1
<i>rmlD</i>	297	dTDP-4-dehydrorhamnose reductase	43.4 (56.2)	42.8 (54.9)	42.1 (55.2)	AHE57213.1
<b>Assembly of the repeat units</b>						
<i>ssB</i>	468	Glucosyl-isoprenylphosphate transferase	41.6 (57.4)	43.4 (69.6)	44.1 (60.1)	AHE57216.1
<i>ssQ</i>	327	Rhamnosyl transferase	22.6 (34.7)	12.5 (17.1)	21.1 (36.6)	AHE56442.1
<i>ssH</i>	428	Glycosyl transferase	N	N	—	AHE52671.1
<i>ssT</i>	289	Glycosyl transferase	N	N	—	AHE56448.1
<i>ssX</i>	652	Acyltransferase	N	N	—	AHE56441.1
<i>ssY</i>	352	Acyltransferase	N	N	—	AHE56447.1
<i>ssZ</i>	372	Acyltransferase	N	N	—	AHE57208.1
<i>ssI</i>	521	Alginate O-acetylation protein	N	N	—	AHE57204.1
<b>Polymerization and export of the repeat units</b>						
<i>ssC-N</i>	1-489/751	Export protein (chain length determinant)	15.6 (29.9)	17.6 (32.3)	16.6 (30.9)	AHE52342.1
<i>ssC-C</i>	490-751/751	Export protein (tyrosine kinase)	20.1 (34.1)	22.8 (36.9)	22.1 (35.9)	AHE52342.1
<i>ssD</i>	199	Polysaccharide export protein	19.5 (26.9)	17.7 (27.8)	17.1 (25.1)	AHE52343.1
<i>ssG</i>	459	Polysaccharide polymerase	17.5 (27.3)	21.9 (34.7)	20.2 (33.6)	AHE52344.1
<i>ssS</i>	482	Polysaccharide biosynthesis protein flippase	14.7 (29.0)	18.7 (33.1)	16.7 (28.5)	AHE57210.1

**Table 1. Identity or similarity analysis between amino acids used for sanxan biosynthesis and the proteins responsible for gellan, welan, and diutan, respectively.** N: not detected; “—” there is no corresponding genome database.

synthesis<sup>32</sup>. However, this connection type ( $\beta$ -D-GlcA  $\leftarrow$   $\beta$ -D-Glc) does not exist in the structure of sanxan (Supplementary Fig. S1). The C2-epimer of  $\beta$ -D-glucose is  $\beta$ -D-mannose, thus *NX02\_24200* might be *ssT*, whose product adds  $\beta$ -D-mannose to  $\beta$ -D-glucuronic acid to form a new type of intermediate ( $\beta$ -D-GlcA  $\leftarrow$   $\beta$ -D-Man). Mutation of *ssT* and *ssQ* was lethal if sanxan synthesis had been initiated on the lipid carrier, which was similar to the knockout of gene *gumB/C/E/M/I/J* in *Xanthomonas campestris*<sup>12</sup>.

Sanxan production was blocked significantly by inactivation of *NX02\_04645*. The phenotype of *NX02* ( $\Delta$ 04645) was *Ss*<sup>-</sup> and could only be recovered by complementation with plasmid pBRR04645, but not by any other glycosyltransferase in the *gel*, *wel* or *sps* clusters. It is speculated that a glycosyltransferase, *SsH*, encoded by *NX02\_04645*, catalyzes the connection of  $\beta$ -D-glucuronic acid to  $\alpha$ -L-Rhamnose, a connection that does not exist in other sphingans. The gene loci of these four GTs are shown in Fig. 2. Therefore, according to the repeat unit of sanxan, four monosaccharides are transferred to the lipid carrier by *SsB*, *SsQ*, *SsH*, and *SsT* in that order.

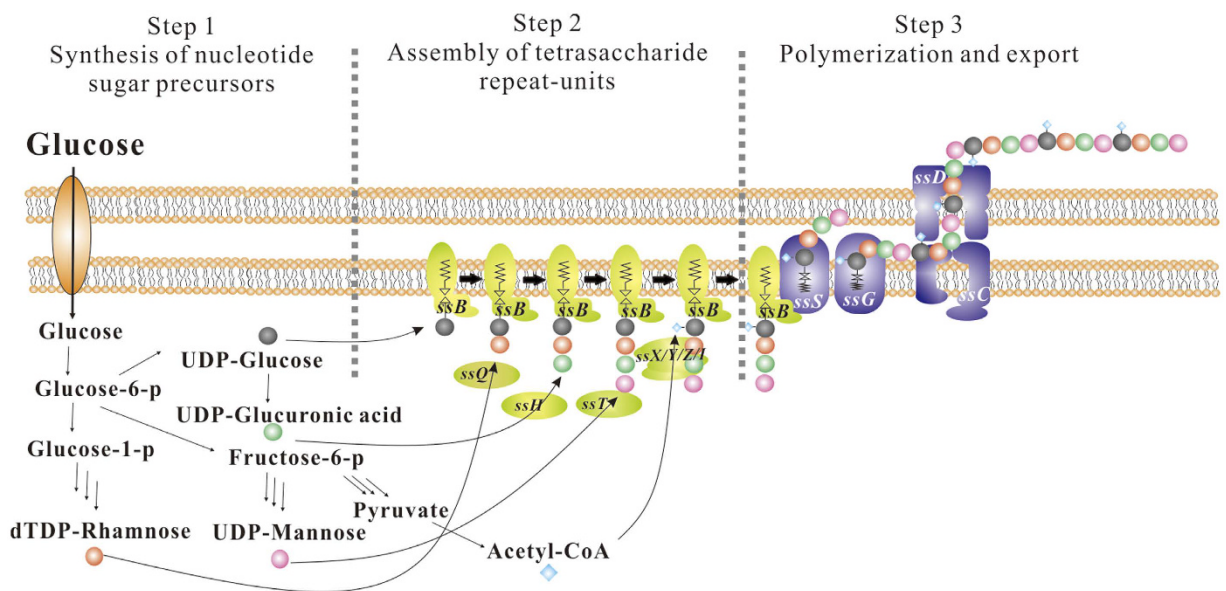
Four ATs genes were found near the GTs and other related genes: *NX02\_24165*, *NX02\_24195*, *NX02\_28130*, and *NX02\_28110*. The phenotypes of the mutants of these genes were all *Ss*<sup>-</sup> and they were complemented by their respective plasmids (Fig. 4). These four ATs were all membrane proteins possessing at least 10 transmembrane domains. These enzymes add acyl groups to the integrated tetrasaccharide repeat unit to prepare for subsequent polymerization and export process. However, the detailed mechanism for the addition of the acyl to the repeat unit is unclear. The schematic diagram of the process of sanxan biosynthesis is shown in Fig. 5.

**Genes and enzymes involved in the polymerization and export of repeat units.** The polymerization and export of sanxan repeat units is a *Wzx/Wzy*-dependent process. No flippase (termed *SsS* in *S. sanxanigenens*) was found in the genome (because of lower identity) by the programme *tblastn* in Bioedit software based on the *gel*, *wel*, and *dps* clusters. While the protein encoded by *NX02\_28140* (K03328, FPKM 180.78) was a polysaccharide transporter, it showed highest (40%) identity with *RfbX*, which is involved in the export of the O-antigen and lower identity with *Gels* (14.7%), *Wels* (18.7%) and *DpsS* (16.7%), respectively (Table 1). The mutated strain could only be obtained in an *ssB*-deficient strain, while it was lethal in the wild-type strain. The phenotype of strain ( $\Delta$ *ssB*,  $\Delta$ 28140) could be complemented by recombinant plasmid pBRRs28140 (Fig. 4). Thus, the protein encoded by *NX02\_28140* is *SsS*.

Two operons encoding polysaccharide co-polymerases, a tyrosine phosphatase, and an outer membrane auxiliary protein were found in the genome, they were *NX02\_02920-02935* and *NX02\_16695-16710* (Fig. 2). *NX02\_16695/02930*, *NX02\_16705/02925*, and *NX02\_16710/02925* were similar to *wza*, *wzc*, and *wzb*, respectively. Notably, *NX02\_02925* and *NX02\_16705/16710* were different. *NX02\_16705/16710* are similar to *gelC/gelE*

Strains	$\Delta ssB$	$\Delta ssB$ $\Delta ssQ$	$\Delta ssM$	$\Delta ssB$ $\Delta ssN$	$\Delta ssX$	$\Delta ssY$	$\Delta ssZ$	$\Delta ssI$	$\Delta ssC$	$\Delta ssD$	$\Delta ssG$	$\Delta ssB$ $\Delta ssS$
	Glycosyltransferase				Acyltransferase				polymerization/export			
Phenotype of mutants												
PCR validation	MWT $\Delta$	MWT $\Delta$	MWT $\Delta$	MWT $\Delta$	MWT $\Delta$	MWT $\Delta$	MWT $\Delta$	MWT $\Delta$	MWT $\Delta$	MWT $\Delta$	MWT $\Delta$	MWT $\Delta$
Complementary plasmids	pBBRB	pBBRSr gelQ	pBBRM	pBBRSr gelL	pBBRX	pBBRY	pBBRZ	pBBRI	pBBRC	pBBRD	pBBRG	pBBRSr S
Phenotype of complementary strains												
The observation of strains	wild-type strain 				The TEM observation of wild-type strain 				The TEM observation of Ss deficient strain 			

**Figure 4.** The phenotypes and genotypes of different mutants and complementary strains, together with transmission electron microscopy observations of wild-type and sanxan-deficient strains.



**Figure 5.** Schematic diagram of sanxan biosynthesis, modified according to Fialho *et al.*<sup>10</sup>. The first step is the synthesis of nucleotide sugar precursors from glucose. The second step is the assembly of the tetrasaccharide repeat unit by the sequential activity of SsB, SsQ, SsH, and SsT glycosyltransferases, and SsX, SsY, SsZ, and SsI acyltransferases. The third step, comprising polymerization and export of the final product, is accomplished by SsS, SsG, SsC, and SsD.

in *S. elodea* ATCC31461, which exhibited a common genetic organization in the *Sphingomonas* genus and is homologous to *NX02\_02925*, the product of which comprises only one polypeptide instead of two independent polypeptides like GelC/GelE<sup>31</sup>. The comparison of their FPKM (Fig. 3) showed higher values for operon *NX02\_02920~02935*. The other operon might be silent or suppressed by certain factors, which was confirmed when its inactivation did not affect the sanxan yield, the viscosity of the fermentation broth, and the product composition. In addition, *NX02\_24775*, located in a gene island, was also analogous to polysaccharide export protein Wza. Its mutant strain had an Ss<sup>+</sup> phenotype. The markerless deletion of operon *NX02\_02920~02935*, or of each gene, affected the phenotype (Ss<sup>-</sup>) significantly, reducing the yield of sanxan (Fig. 4). The polypeptide named SsC, encoded by *NX02\_02925*, was identified as the autophosphorylating tyrosine kinase involved in polysaccharide chain length determination<sup>31</sup>. The Ss<sup>-</sup> phenotype of strain *NX02* ( $\Delta$ *ssC*) could not be completely recovered by the transformation with plasmid pBBRssC and the yield of sanxan was slightly decreased. SsC comprises 751 amino acids, of which M<sub>1</sub> to R<sub>489</sub> from the N-terminus are homologous to GelC (15.6%), WelC (17.6%), and DpsC (16.6%), and E<sub>490</sub> to G<sub>751</sub> from the C-terminus are similar to GelE (20.1%), WelE (22.8%), and DpsE (22.1%) (Table 1); these two segments are homologous to the activator domain and the kinase domain of SsC, respectively. SsC was predicted to have two transmembrane  $\alpha$ -helices, TM1 and TM2, located at W<sub>54</sub> to T<sub>76</sub> and V<sub>469</sub> to A<sub>488</sub>. *NX02\_02930*, named as *ssD*, is homologous to *spnD*, the product of which is an OMA protein homologue that is responsible for the export of sanxan chains. *NX02* ( $\Delta$ *ssD*, Ss<sup>-</sup>) could be complemented by plasmid pBBRssD to an Ss<sup>+</sup> phenotype. SsD was not predicted to have a transmembrane helix, like Gumb; however, GelD, WelD, and DpsD all have one helix in their N-terminus. Although three copies of *wza* and two copies of *wzc* homologous genes were found in the genome, only one gene was responsible for sanxan biosynthesis.

The polymerase related to sanxan biosynthesis was named SsG. *NX02\_02935*, annotated to encode an O-antigen polymerase, was identified by browsing the whole genome. Its deletion reduced the production of sanxan drastically, thus *NX02\_02935* was gene *ssG*. However, multicopy expression of *ssG* in *NX02* ( $\Delta$ *ssG*) did not recover the Ss<sup>-</sup> phenotype after transformation with pBBRssG (Fig. 4). This might be because a balanced expression level of *ssC* and *ssG* is necessary to assemble the membrane protein complex correctly. In addition, the Ss<sup>-</sup> phenotype of strains *NX02* ( $\Delta$ *ssC*), *NX02* ( $\Delta$ *ssD*), and *NX02* ( $\Delta$ *ssG*) could not be complemented by plasmids pBBRgelC/E (*welC/E*), pBBRgelD (*welD*), and pBBRgelG (*welG*), respectively. Thus, the polysaccharide biosynthesis process might show catalytic specificity for the polymerization and export of the repeat units.

**Regulatory gene.** A multi-sensor hybrid histidine kinase SsA, encoded by *NX02\_06855*, is homologous with GelA (60.5% identity) from *S. elodea*. It contains 797 amino acids and two transmembrane helices in the N terminus: G<sub>22</sub> to G<sub>44</sub> and G<sub>49</sub> to F<sub>66</sub>. Our knockout experiments showed that mutation of *ssA* did not affect the yield of sanxan. In addition, the expression levels of genes related to sanxan synthesis between *NX02* ( $\Delta$ *ssA*) and the wild-type were similar or only slightly altered (Supplementary Fig. S4). Thus, the deletion of *ssA* did not affect the expression levels of related genes. Therefore, the positive regulator of gellan, welan and other sphingans synthesis could not regulate sanxan production.

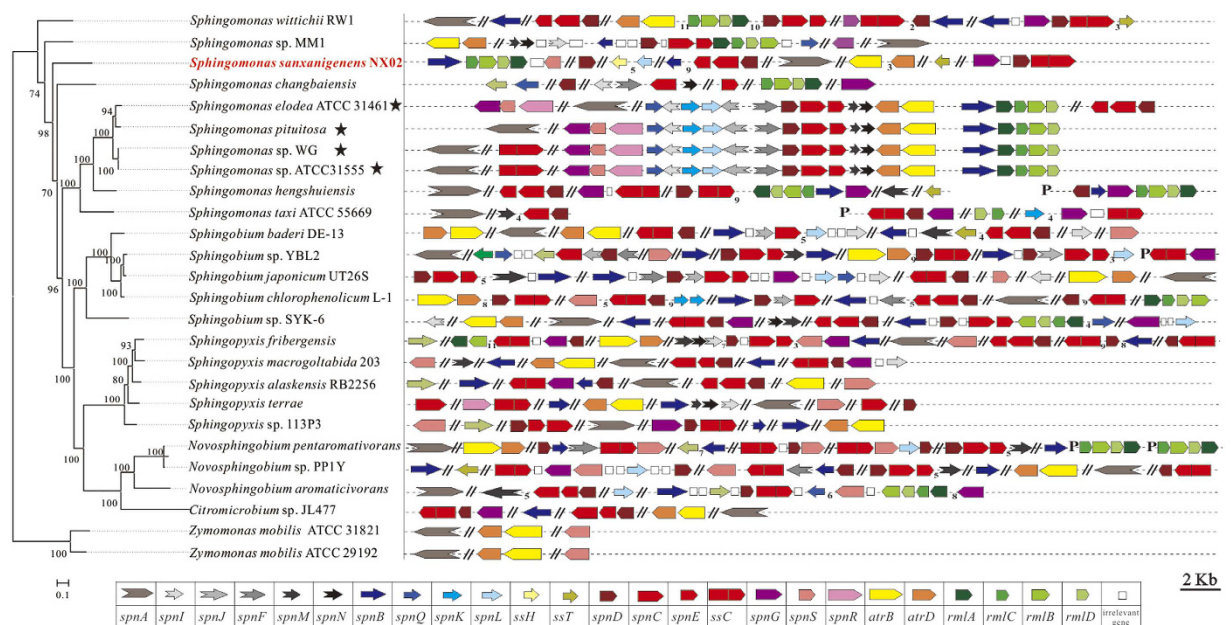
**Sanxan is a capsular polysaccharide.** The sphingans, such as gellan and welan, are structurally related EPSs secreted by a group of the genus *Sphingomonas*<sup>9</sup>. A lyase, SpnR, found in the *spn* cluster, released the polysaccharide from the outer membrane into extracellular environment. In addition, the deletion of *dpsM* and/or *dpsN* led to more easy removal of the polysaccharide from the cells<sup>50</sup>. However, *spnR*, *dpsM*, and *dpsN* were all not found in *NX02* genome. The surfaces of bacteria cells of *S. elodea* ATCC31461 and *Sphingomonas* sp. ATCC 31555 were both smooth in YEME medium, while in *NX02*, sanxan was spread over the cell surface (Supplementary Fig. S5a,b,c). When plasmid pBBRgelR was overexpressed in *NX02*, sanxan was released from the *NX02* cell surface and the capsular-free cells tended to gather together (Supplementary Fig. S5d). These results suggested that sanxan is a capsular polysaccharide and that such a lyase does not exist in *NX02* genome.

**Comparisons of the biosynthetic pathways of sanxan, sphingans, and xanthan.** All gene loci related to sanxan biosynthesis are shown in Supplementary Fig. S6, and are scattered over the whole genome. By contrast, the genes responsible for the assembly, polymerization, and export in all sphingans were clustered, with a uniform arrangement<sup>6,11</sup>, although they are transcribed by several promoters. Twelve *gum* genes form an operon under the control of a single promoter. Four *rml* genes were clustered and arranged as the sequence *rmlC-B-D-A*, showed the same order with *S. wittichii*, while those in sphingan-producing strain were all *rmlA-C-B-D*<sup>11</sup>. In addition, the *spnI*, *spnJ*, *spnF*, *spnM*, *spnN*, and *spnR* genes were not found in the *NX02* genome, which suggest these six genes were not essential for the biosynthesis of sanxan during the long-term evolution. Two submits of protein SsC are homologous with SpnC and SpnE. Four ATs were also found in *NX02* genome. In brief, the biosynthetic pathway of sanxan was more like a “patchwork” of dispersed gene elements from different locations. Thus, because of the obviously different pathways and the low similarity of the related genes, the structure of sanxan is distinct.

### The evolutionary analysis of the arrangement of genes related to sanxan biosynthesis.

Compared with sphingans, the capsular polymer sanxan possesses a specific structure, specific properties<sup>34</sup>, and scattered biosynthetic genes. In addition, *NX02* is phylogenetically distant from sphingan producing strains, and the homology between the *ss* genes (all genes related to sanxan biosynthesis) and *spn* genes was very low. Furthermore, many unstable genetic elements exist in the *NX02* genome. To obtain clues to the evolutionary process of the arrangement of genes related to sanxan biosynthesis, we analysed all those genes in *Sphingomonadaceae* strains with completely sequenced genomes, based on gene annotation and homology alignment against *ss* and *gel* genes. The genes related to the assembly, polymerization and export and its distribution in 26 genomes, include 22 completed sequenced genomes and four genomes of sphingan producing strains (*Sphingomonas* sp. ATCC 31555, *Sphingomonas* sp. WG, *S. elodea* ATCC 31461 and *Sphingomonas pituitosa*) are shown in Fig. 6. The arrangement





**Figure 6. Genes related to assembly, polymerization and export, and their distribution in 26 genomes.** The phylogenetic relationships of *S. sanxanigenens* NX02 with 25 sequenced strains of the family *Sphingomonadaceae* (including 21 other sequenced genomes and four sphingane-producing genomes) were constructed using the Neighbor-Joining method based on all core genes. Related genes were predicted by gene annotation and homology alignment against *ss* and *gel* genes. “★” indicates the sphingane-producing strains.

of genes in the four sphingane-producing strains was approximately consistent, except in one case where the location of *gelG/S/R* was distant from the main area. Although 21 strains were reported as non-sphingane producing strains, most genes existed in a dispersed form in their genomes, which was similar to NX02. For example, *S. wittichii* RW1 possessed most genes except for one *GT*; and *Sphingomonas* sp. MM1, *S. japonicum* UT26S, *Sphingobium* sp. SYK-6, and *N. aromativorans* retained many traces of the *spn* cluster. Genes *spnD/C/E* were always clustered and were present as multicopies, which suggests their evolutionary diversity. The order of *rml* genes in NX02 was same as that in sphingane-free strains.

Therefore, we deduced that in the progenitor of NX02, *ss* genes were scattered over the genome and were incomplete like other non-sphingane producing strains, and the phenotype of progenitor NX02 might have been  $Ss^-$ . Sanxan could be produced when some adaptive mutations and HGTs happened during long term evolution. The large number of GIs, ISs, transposases, and prophages would facilitate this process of evolution. The *ssH* gene is located between a GI and an IS (Fig. 2), thus it might be an exogenous gene obtained by HGT. Judging by their specificity, adaptive mutations might have occurred in the nucleotide sequence of the *ssC/D/G/S* genes. The evolution of the synthetic pathway for sanxan was independent and very different from that of sphingane. Consequently, the structure of sanxan is obviously different from the sphinganes, and the common positive regulator *SsA* of sphinganes could not regulate sanxan production. Based on the above, we concluded that *S. sanxanigenens* NX02 was a natural polysaccharide producing strain that evolved over a long-time.

**A hypothetical evolutionary model related to the polysaccharide biosynthetic cluster.** Genes responsible for polysaccharide biosynthesis were always clustered<sup>5–8,11</sup>. Research into the evolution of their pathways is difficult because the phenotype of a strain would change to  $EPS^-$  when the cluster was destroyed during the long period evolution, and these mutants with the  $EPS^-$  phenotype would not arouse researchers' attention. Itoh *et al.*<sup>18</sup> pointed out that the shuffling of a genome structure was virtually neutral over long-term evolution and that gene order in operons was unstable. This evolutionary process also has been demonstrated by analysis of the *cps* (capsular polysaccharide synthesis) gene clusters within *Klebsiella* spp., in which many shuffling phenomena such as lateral gene transfer, truncation, and transposition, were observed<sup>5</sup>. Therefore, genes in a cluster are also not constant; events such as translocation, deletion, and insertion happen frequently (Fig. 7I). Subsequently, cluster 1 changed to a, b, c, and d forms or with a different gene arrangement, as described in Fig. 7, and polysaccharides could not be produced after complete destruction of the cluster (Fig. 7e). However, a metabolic pathway should be regulated for the effective use of energy, with only the related genes being organized into operons or clusters<sup>17</sup>, and this gene cluster could also promote the lateral transfer of the phenotype<sup>1,13,14</sup> (Fig. 7II). In addition to these two models of cluster formation, other models have been proposed for the formation of operons, a more structured instance of cluster, for example the Natal Model<sup>51</sup>, Fisher Model<sup>2</sup>, and Co-regulation Model<sup>16</sup>. Furthermore, an evolutionary model for the origin and evolution of proteobacterial histidine biosynthetic operons described a piecemeal building process from single genes to one operon<sup>52</sup>. After long-term evolution, a new cluster 2 would appear that did not include the non-essential gene *D* or essential genes with different arrangement

Strain or plasmid	Genotype or phenotype	Source or reference
<b>Strains</b>		
<i>E. coli</i> S17	<i>RecA thi pro hsdR<sup>-</sup> M<sup>+</sup> RP4, Sm<sup>R</sup> Amp<sup>R</sup> Kan<sup>R</sup></i>	This work
<i>S. sanxanigenens</i> NX02	Wild-type strain, Ss <sup>+</sup> , Cm <sup>R</sup>	This work
<i>S. elodea</i> ATCC31461	Wild-type strain, Gel <sup>+</sup>	This work
<i>Sphingomonas</i> sp. ATCC31555	Wild-type strain, Wel <sup>+</sup>	This work
<b>Plasmids</b>		
pLO3	4937-bp suicide vector, tet <sup>R</sup>	53
pBBR1MCS-2	5144-bp broad host range vector, kan <sup>R</sup>	54
pLO3ssn	pLO3 derivative carrying upstream and downstream fragment of <i>ssn</i> , “ <i>ssn</i> ” refers to all gene related to Ss biosynthesis	This work
pBBRssB	pBBR1MCS-2 derivative expressing <i>ssB</i>	This work
pBBRssn	pBBR1MCS-2 derivative expressing <i>ssn</i> , “ <i>ssn</i> ” refers to all related genes	This work
pBBRsrssn	pBBR1MCS-2 derivative expressing <i>ssB</i> and <i>ssn</i> simultaneously	This work
pBBRrgeln	pBBR1MCS-2 derivative expressing <i>ssB</i> and <i>geln</i> simultaneously, “ <i>geln</i> ” refers to <i>gelQ</i> , <i>gelL</i>	This work
pBBRgeln	pBBR1MCS-2 derivative expressing <i>geln</i> , “ <i>geln</i> ” refers to <i>gelD/gelCE/gelS/gelG</i>	This work
pBBRweln	pBBR1MCS-2 derivative expressing <i>weln</i> , “ <i>weln</i> ” refers to <i>welD/welCE/welS/welG</i>	This work

**Table 2. Bacterial strains and plasmids used in this study.**

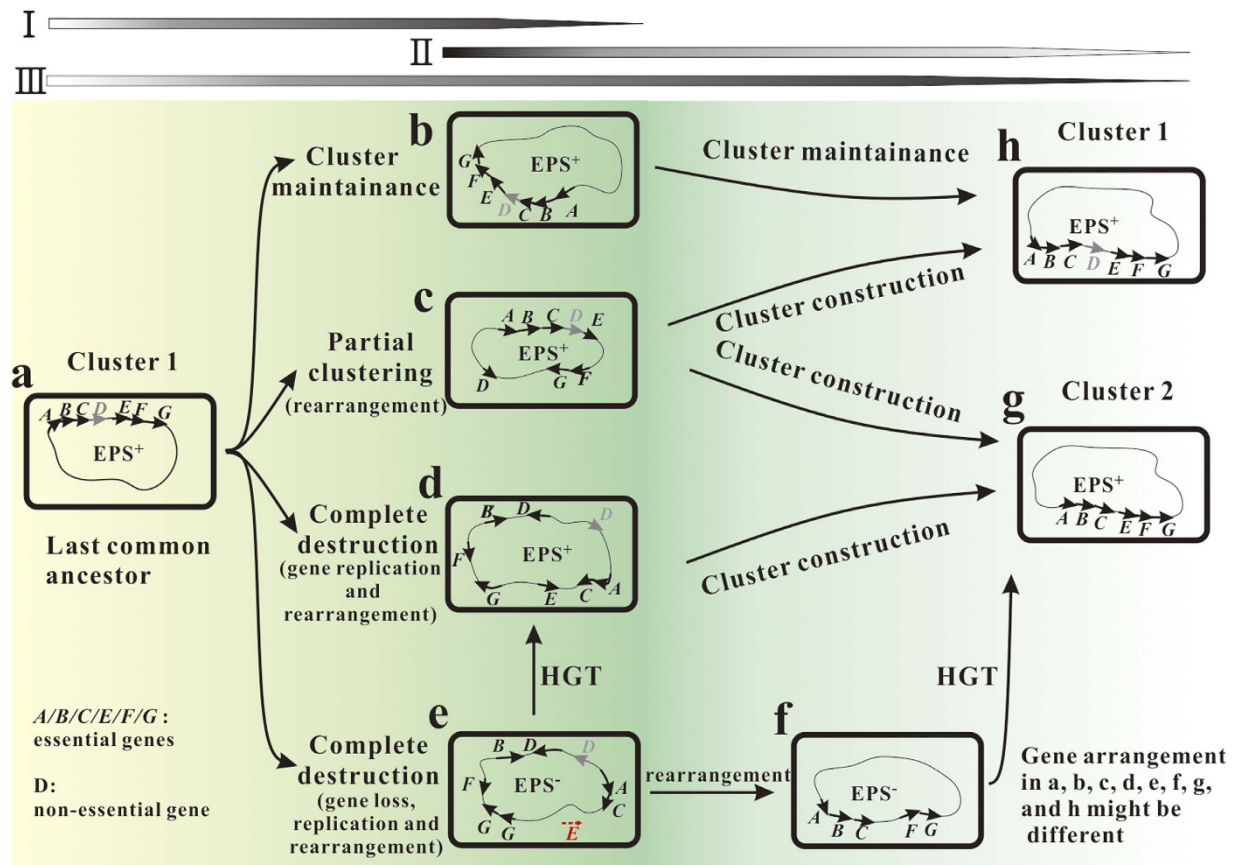
(Fig. 7g). Therefore, the evolutionary process of the biosynthetic pathway for microbial polysaccharide might be proposed as a lengthy cyclic process: cluster 1 → scatter → cluster 2 (Fig. 7III). In this process, genes in cluster 1 and cluster 2 were always not identical. The biosynthetic pathway of sanxan proved the existence of the scatter process. This process would lead the appearance of many new species.

Based on *Sphingomonadaceae*, the gene cluster of the sphingan-producing strains might be the last common ancestor (Fig. 7). With increasing time, some genes in this cluster were translocated, deleted, and the cluster might have been broken or lost. Strong evidence for this hypothesis is provided by the existence of partial non-essential *spnI/J/F/M/N/R* genes in most sphingan-free strains, such as *Sphingomonas* sp. MM1 and *S. japonicum* UT26S, and the translocation of the *gelG/S/R* genes in the *gel* cluster (Supplementary Fig. S6). The phenotype of EPS<sup>+</sup> was lost during this process, and many new strains appeared. NX02 might have undergone this evolutionary journey, and losing the six *spnI/J/F/M/N/R* genes in the process. Subsequently, to defend against extreme environments, NX02 acquired some genes by HGT or adaptive mutations to produce a capsular polysaccharide that is different from the sphingans. Thus, the organization of *ss* genes demonstrated that clustering is not essential for polysaccharide production. However, as a more efficient form, a new cluster or operon will appear after long-term evolution. The new cluster might not contain the *spnI/J/F/M/N/R* genes, and will be stable under the pressure of severe environments, or might be destroyed because of the unstable genome structure. The putative evolution process of NX02 could be described as: a → e → d → g (Fig. 7). It is likely that after long-term evolution, another polysaccharide-producing strain like NX02 in the *Sphingomonadaceae* will appear.

## Methods

**Strains, plasmids, media, and culture conditions.** The bacterial strains and plasmids used in this work are listed in Table 2. *Escherichia coli* strains were grown in Luria Bertani (LB) medium. *Sphingomonas Sanxanigenens* NX02 was cultured at 30 °C on NKG (NK medium with 1.5% glucose; NK: 0.5% peptone, 0.3% beef powder, 0.1% yeast extract and 1.5% agar 15.0 g, pH 7.0), NKS (NK medium with 8% sucrose), and YEME medium (0.25% yeast extract, 0.025% malt extract) that was developed to reduce sphingan production and improve cell suspension<sup>28</sup>. Antibiotics were used at the following concentrations (μg/mL): tetracycline (Tc; 10), kanamycin (Km; 25), and Chloramphenicol (Cm; 25). For sanxan fermentation, cells were grown in medium consisting of: 4% glucose, 0.02% yeast extract, 0.12% K<sub>2</sub>HPO<sub>4</sub>, 0.2% NaNO<sub>3</sub>, 0.1% CaCO<sub>3</sub>, 0.0005% FeSO<sub>4</sub>, 0.04% NaCl, and 0.05% MgSO<sub>4</sub> (pH 7.5)<sup>55</sup>. Peptone, beef powder, yeast extract, agar, and other chemicals were purchased from Dingguo Limited (Tianjin, China).

**Genome sequencing and analysis of ORFs from *S. sanxanigenens* NX02.** The complete genome sequence of *S. sanxanigenens* NX02 has been deposited in GenBank under the accession nos. CP006644 and CP011450. Whole genome sequencing was performed using the Illumina HiSeq 2000 and Pacific RSII platforms. The genome was assembled using 1.2 GB Illumina paired-end reads, 1.7 GB Illumina mate-paired reads, and 53.8 MB PacBio reads. Sequence quality assessment and assembly were performed with a quality of <1 error in 100,000 bases using PHRAP and Consed. Error correction of the PacBio reads was performed using the Illumina reads. Genes were predicted using Glimmer3<sup>56</sup> and tRNAscan-SE<sup>57</sup>, and annotated by searching against the nr protein database of GenBank using blastp, with E values less than 1.00<sup>e-5</sup>. A Neighbour-Joining phylogenetic



**Figure 7. Hypothetical scheme for the lengthy cyclic evolutionary process of the biosynthetic pathway for microbial polysaccharides.** There are many possible forms (a–g) of the evolutionary process (I–III). “A~G” stands for different genes related to polysaccharide biosynthesis.

tree was constructed in Mega<sup>58</sup>. Genomic islands were predicted using IslandViewer, which integrated the IslandPath-DIMOB and SIGI-HMM algorithms<sup>59</sup>. The ISs were identified and classified using the ISfinder database<sup>60</sup>. Percent identities or similarities between amino acid sequences were calculated using the online programme EMBOSS Needle, ([http://www.ebi.ac.uk/Tools/psa/emboss\\_needle/](http://www.ebi.ac.uk/Tools/psa/emboss_needle/)). The prediction of transmembrane helices in proteins was performed using the TMHMM Server v.2.0 (<http://www.cbs.dtu.dk/services/TMHMM/>). Genes associated with certain pathways were analysed using the Kyoto Encyclopedia of Genes and Genomes (KEGG) database (<http://www.genome.jp/kegg/>). The glycosyltransferases of *S. sanxanigenens* NX02 were analysed using the Carbohydrate-Active Enzymes database (CAZY; <http://www.cazy.org/>) and the NCBI database.

**Markerless gene knockout and complementation.** Genes were inactivated by double-crossover homologous recombination. The upstream and downstream flanking sequences (approximate 1.5 kb) of genes were spliced using overlap extension PCR. The primers used to amplify the flanking fragments of the target genes contained SacI, XbaI, or PacI restriction sites; all primers are shown in Supplementary Table S4. PCR products of the respective genes were digested with the appropriate restriction enzymes, ligated into a suicide vector, pLO3<sup>53</sup> (Table 2), and used to transform *E. coli* S17. The respective recombinant plasmids were transferred to *S. sanxanigenens* NX02 wild-type strain or NX02 ( $\Delta$  *ssB*) strain using biparental filter mating at 30 °C for 12 h on NKG medium without antibiotics. The single crossover mutants were selected on NKG medium containing 10 µg/mL Tc and 25 µg/mL Cm. The knockout mutants were then isolated on NKS medium with 25 µg/mL Cm, followed by PCR screening using the verification primers. The large fragments were deleted using the same procedures.

To identify the function of a deleted gene, complementation tests were performed as follows. The targeted genes related to sanxan synthesis, or specific genes from *S. elodea* ATCC31461 or *Sphingomonas* sp. ATCC 31555, were amplified and ligated into the broad host range expression vector pBBR1MCS-2<sup>54</sup> or pBBRssB (Table 2), and the recombinant vectors were transferred into the respective mutant strains of *S. sanxanigenens* NX02 using biparental conjugation. Primers used to construct expression vectors are shown in Supplementary Table S4. The recombinant expression vectors were verified by PCR screening and DNA sequencing, and the recombinant strains harbouring the plasmids were selected using PCR. The genome of *S. sanxanigenens* NX02 was extracted using an AxyPreP™ Bacterial Genomic DNA Miniprep Kit (Axygen, Hangzhou, China). Plasmid DNA was purified from *E. coli* using the Axyprep™ Plasmid Miniprep Kit (Axygen). PCR products were purified using a DNA Gel Extraction Kit (Axygen) and a PCR Cleanup Kit (Axygen).

**RNA isolation and chain specific transcriptome sequencing.** Large amounts of sanxan accumulated around the cells on NK medium. The crude total DNA-free RNA of *S. sanxanigenens* NX02 was extracted using the RNAiso Plus (Takara, Dalian, China) and RNAPrep Pure Cell/Bacteria Kit (Tiangen, China) when the strain reached cultured at logarithmic phase in NK medium. Total RNA quality was assessed using a gel electrophoresis BioDrop Cuvette (BioDrop, United Kingdom). rRNA was depleted from the total RNA using a Ribo-Zero Magnetic kit (Epicentre, Madison, WI, USA). Chain specific transcriptome sequencing of the double-stranded cDNA was performed following the Illumina workflow on a HiSeq 2500 (Illumina) using the Truseq PE Cluster Kit v3-cBot-HS (Illumina) and the cBot instrument (Illumina)<sup>61</sup>. A total of 17,549,424 reads were generated that resulted in 285-fold sequencing coverage. The sequence quality satisfied the criterion of <3 error in 10,000 bases. The number of fragments per kilobase of exon per million mapped reads (FPKM) was calculated to measure expression levels of the genes by RSEM (RNA-Seq by Expectation-Maximization, <http://deweylab.biostat.wisc.edu/rsem/>)<sup>62,63</sup>. Operon identification was also performed using chain specific transcriptome sequencing<sup>64</sup>. If multiple genes share the same transcriptional start site and termination site after expanded sweep, these genes will belong to an operon<sup>64</sup>.

**cDNA synthesis, and qRT-PCR.** Total RNA (1.5 µg) was reverse-transcribed using a Quantscript RT Kit (Tiangen, China), according to the manufacturer's protocol. The relative expression analysis of genes related to sanxan biosynthesis in different strains was performed using the quantitative RT-PCR with a MyiQ™ two-colour real-time PCR detection system (BIO-RAD laboratories) with the Bestar® SybrGreen qPCR mastermix (DBI, Bioscience Inc., Germany). The primers were designed using the OLIGO software and the length of amplicons was between 100 and 200 bp. The primer sequences used in qRT-PCR are listed in Supplementary Table S5. The endogenous reference gene was 16srRNA. Standard deviations were calculated from three PCR replicates and the relative abundance of the genes was determined using the comparative Ct method.

**Analysis of fermentation broth.** The extraction of sanxan from the fermentation was performed according to a previously published method<sup>35</sup>. The viscosity of the sanxan solution was measured using a Brookfield viscometer DV\_II + (USA) equipped with a no. 64 spindle at a shear rate of 60 rev/min.

**Electron microscopy.** Strains with different phenotypes or genotypes were prepared for transmission electron microscopy (TEM; Hitachi, Tokyo, Japan). NX02, *Sphingomonas elodea* ATCC31461 and *Sphingomonas* sp. ATCC31555 strains were cultured in YEME medium at 30 °C for 18 h. Then strains were collected, washed with phosphate buffer twice to remove impurities, and 1 µl of cell suspension at an appropriate concentration was dropped onto Holey carbon Film and observed directly<sup>65</sup>.

## References

- Price, M. N., Arkin, A. P. & Alm, E. J. The life-cycle of operons. *PLoS Genet.* **2**, e96 (2006).
- Martin, F. J. & McInerney, J. O. Recurring cluster and operon assembly for phenylacetate degradation genes. *BMC Evol. Biol.* **9**, 1–11 (2009).
- Wong, S. & Wolfe, K. H. Birth of a metabolic gene cluster in yeast by adaptive gene relocation. *Nat. Genet.* **37**, 777–782 (2005).
- Dandekar, T., Snel, B., Huynen, M. & Bork, P. Conservation of gene order: a fingerprint of proteins that physically interact. *Trends Biochem. Sci.* **23**, 324–328 (1998).
- Pan, Y. J. *et al.* Genetic analysis of capsular polysaccharide synthesis gene clusters in 79 capsular types of *Klebsiella* spp. *Sci. Rep.* **5**, 15573 (2015).
- Schmid, J., Sieber, V. & Rehm, B. Bacterial exopolysaccharides: biosynthesis pathways and engineering strategies. *Front. Microbiol.* **6**, 496 (2015).
- Ji, K. *et al.* Bacterial cellulose synthesis mechanism of facultative anaerobe *Enterobacter* sp. FY-07. *Sci. Rep.* **6**, 21863 (2016).
- Rehm, B. H. A. & Valla, S. Bacterial alginates: biosynthesis and applications. *Appl. Microbiol. Biotechnol.* **48**, 281–288 (1997).
- Glucksmann, M. A., Reuber, T. L. & Walker, G. C. Genes needed for the modification, polymerization, export, and processing of succinoglycan by *Rhizobium meliloti*: a model for succinoglycan biosynthesis. *J. Bacteriol.* **175**, 7045–7055 (1993).
- Fialho, A. M. *et al.* Occurrence, production, and applications of gellan: current state and perspectives. *Appl. Microbiol. Biotechnol.* **79**, 889–900 (2008).
- Schmid, J., Sperl, N. & Sieber, V. A. Comparison of genes involved in sphingane biosynthesis brought up to date. *Appl. Microbiol. Biotechnol.* **98**, 7719–7733 (2014).
- Becker, A., Katzen, F., Pühler, A. & Ielpi, L. Xanthan gum biosynthesis and application: a biochemical/genetic perspective. *Appl. Microbiol. Biotechnol.* **50**, 145–152 (1998).
- Lawrence, J. G. & Roth, J. R. Selfish operons: horizontal transfer may drive the evolution of gene clusters. *Genetics* **143**, 1843–1860 (1996).
- Lawrence, J. Selfish operons: the evolutionary impact of gene clustering in prokaryotes and eukaryotes. *Curr. Opin. Genet. Dev.* **9**, 642–648 (1999).
- Ballouz, S., Francis, A. R., Lan, R. & Tanaka, M. M. Conditions for the evolution of gene clusters in bacterial genomes. *PLoS Comput. Biol.* **6**, e1000672 (2010).
- Price, M. N., Huang, K. H., Arkin, A. P. & Alm, E. J. Operon formation is driven by co-regulation and not by horizontal gene transfer. *Genome Res.* **15**, 809–819 (2005).
- Svetic, R. E., MacCluer, C. R., Buckley, C. O., Smythe, K. L. & Jackson, J. H. A metabolic force for gene clustering. *B. Math. Biol.* **66**, 559–581 (2004).
- Itoh, T., Takemoto, K., Mori, H. & Gojbori, T. Evolutionary instability of operon structures disclosed by sequence comparisons of complete microbial genomes. *Mol. Biol. Evol.* **16**, 332–346 (1999).
- Yabuuchi, E. *et al.* Proposal of *Sphingomonas wittichii* sp. nov. for strain RW1T, known as a dibenzo-p-dioxin metabolizer. *Int. J. Syst. Evol. Micr.* **51**, 281–292 (2001).
- Takeuchi, M., Hamana, K. & Hiraishi, A. Proposal of the genus *Sphingomonas* sensu stricto and three new genera, *Sphingobium*, *Novosphingobium* and *Sphingopyxis*, on the basis of phylogenetic and chemotaxonomic analyses. *Int. J. Syst. Evol. Micr.* **51**, 1405–1417 (2001).

21. Yabuuchi, E. *et al.* Emendation of the genus *Sphingomonas* Yabuuchi *et al.* 1990 and junior objective synonymy of the species of three genera, *Sphingobium*, *Novosphingobium* and *Sphingopyxis*, in conjunction with *Blastomonas ursincola*. *Int. J. Syst. Evol. Microbiol.* **52**, 1485–1496 (2002).
22. Busse, H. J. *et al.* *Sphingomonas aurantiaca* sp. nov., *Sphingomonas aerolata* sp. nov. and *Sphingomonas faeni* sp. nov., air- and dustborne and Antarctic, orange-pigmented, psychrotolerant bacteria, and emended description of the genus *Sphingomonas*. *Int. J. Syst. Evol. Microbiol.* **53**, 1253–1260 (2003).
23. Whitfield, C. Biosynthesis and assembly of capsular polysaccharides in *Escherichia coli*. *Annu. Rev. Biochem.* **75**, 39–68 (2006).
24. Schmid, J. & Sieber, V. Enzymatic transformations involved in the biosynthesis of microbial exo-polysaccharides based on the assembly of repeat units. *ChemBioChem* **16**, 1141–1147 (2015).
25. Yamazaki, M., Thorne, L., Mikolajczak, M., Armentrout, R. W. & Pollock, T. J. Linkage of genes essential for synthesis of a polysaccharide capsule in *Sphingomonas* strain S88. *J. Bacteriol.* **178**, 2676–2687 (1996).
26. Thorne, L., Mikolajczak, M. J., Armentrout, R. W. & Pollock, T. J. Increasing the yield and viscosity of exopolysaccharides secreted by *Sphingomonas* by augmentation of chromosomal genes with multiple copies of cloned biosynthetic genes. *J. Ind. Microbiol. Biot.* **25**, 49–57 (2000).
27. Sá-Correia, I. *et al.* Gellan gum biosynthesis in *Sphingomonas paucimobilis* ATCC 31461: genes, enzymes and exopolysaccharide production engineering. *J. Ind. Microbiol. Biot.* **29**, 170–176 (2002).
28. Coleman, R. J., Patel, Y. N. & Harding, N. E. Identification and organization of genes for diutan polysaccharide synthesis from *Sphingomonas* sp. ATCC 53159. *J. Ind. Microbiol. Biot.* **35**, 263–274 (2008).
29. Kaur, V., Bera, M. B., Panesar, P. S., Kumar, H. & Kennedy, J. F. Welan gum: microbial production, characterization, and applications. *Int. J. Biol. Macromol.* **65**, 454–461 (2014).
30. Videira, P., Fialho, A., Geremia, R. A., Breton, C. & Sá-correia, I. Biochemical characterization of the  $\beta$ -1, 4-glucuronosyltransferase GelK in the gellan gum-producing strain *Sphingomonas paucimobilis* ATCC 31461. *Biochem. J.* **358**, 457–464 (2001).
31. Moreira, L. M. *et al.* The gellan gum biosynthetic genes *gelC* and *gelE* encode two separate polypeptides homologous to the activator and the kinase domains of tyrosine autokinases. *J. Mol. Microbiol. Biotech.* **8**, 43–57 (2005).
32. Harding, N. E., Patel, Y. N. & Coleman, R. J. Organization of genes required for gellan polysaccharide biosynthesis in *Sphingomonas elodea* ATCC 31461. *J. Ind. Microbiol. Biot.* **31**, 70–82 (2004).
33. Huang, H. D. *et al.* *Sphingomonas sanxanigenens* sp. nov., isolated from soil. *Int. J. Syst. Evol. Microbiol.* **59**, 719–723 (2009).
34. Huang, H. D. *et al.* Structural and physical properties of sanxan polysaccharide from *Sphingomonas sanxanigenens*. *Carbohydr. Polym.* **144**, 410–418 (2016).
35. Wu, M. M. *et al.* The simultaneous production of sphingans Ss and poly (R-3-hydroxybutyrate) in *Sphingomonas sanxanigenens* NX02. *Int. J. Biol. Macromol.* **82**, 361–368 (2016).
36. Hong, H. B., Chang, Y. S., Nam, I. H., Fortnagel, P. & Schmidt, S. Biotransformation of 2, 7-dichloro- and 1, 2, 3, 4-tetrachlorodibenzo-p-dioxin by *Sphingomonas wittichii* RW1. *Appl. Environ. Microbiol.* **68**, 2584–2588 (2002).
37. Copley, S. D. *et al.* The whole genome sequence of *Sphingobium chlorophenicum* L-1: insights into the evolution of the pentachlorophenol degradation pathway. *Genome Biol. Evol.* **4**, 184–198 (2012).
38. Nagata, Y. *et al.* Genomic organization and genomic structural rearrangements of *Sphingobium japonicum* UT26, an archetypal  $\gamma$ -hexachlorocyclohexane-degrading bacterium. *Enzyme Microb. Tech.* **49**, 499–508 (2011).
39. Glaeser, S. P. & Kämpfer, P. *The family sphingomonadaceae in The Prokaryotes* (ed. Rosenberg, E. *et al.*) 641–707 (Springer Berlin Heidelberg, 2014).
40. Gai, Z. *et al.* Genome sequence of *Sphingomonas elodea* ATCC 31461, a highly productive industrial strain of gellan gum. *J. Bacteriol.* **193**, 7015–7016 (2011).
41. Wang, X., Tao, F., Gai, Z., Tang, H. & Xu, P. Genome sequence of the welan gum-producing strain *Sphingomonas* sp. ATCC 31555. *J. Bacteriol.* **194**, 5989–5990 (2012).
42. Li, H. *et al.* Draft Genome Sequence of *Sphingomonas* sp. WG, a Welan Gum-Producing Strain. *Genome Announc.* **4**, e01709–15 (2016).
43. Bahl, M. A., Schultheis, E., Hempel, D. C., Nörtemann, B. & Franco-Lara, E. Recovery and purification of the exopolysaccharide PS-EDIV from *Sphingomonas puitosa* DSM 13101. *Carbohydr. Polym.* **80**, 1037–1041 (2010).
44. Zhang, J. Y., Liu, X. Y. & Liu, S. J. *Sphingomonas changbaiensis* sp. nov., isolated from forest soil. *Int. J. Syst. Evol. Microbiol.* **60**, 790–795 (2010).
45. Chen, I. & Dubnau, D. DNA uptake during bacterial transformation. *Nat. Rev. Microbiol.* **2**, 241–249 (2004).
46. Smeets, L. C. & Kusters, J. G. Natural transformation in *Helicobacter pylori*: DNA transport in an unexpected way. *Trends Microbiol.* **10**, 159–162 (2002).
47. Huang, H. D. *et al.* Cloning, expression and characterization of a phosphoglucosyltransferase/phosphomannosyltransferase from sphingans-producing *Sphingomonas sanxanigenens*. *Biotechnol. Lett.* **35**, 1265–1270 (2013).
48. Marques, A. R., Ferreira, P. B., Sa-Correia, I. & Fialho, A. M. Characterization of the *ugpG* gene encoding a UDP-glucose pyrophosphorylase from the gellan gum producer *Sphingomonas paucimobilis* ATCC 31461. *Mol. Genet. Genomics* **268**, 816–824 (2003).
49. Wu, M. M., Huang, H. D., Li, G. Q., Zhou, J. F. & Ma, T. Biochemical characterization and functional analysis of Udp-glucose dehydrogenase, in the synthesis of biopolymer Ss from *Sphingomonas sanxanigenens* NX02. *Appl. Biochem. Microbiol.* **51**, 27–33 (2015).
50. Harding, N. E., Patel, Y. N. & Coleman, R. J. Targeted gene deletions for polysaccharide slime formers. *U.S. Patent No. 9,422,567*. Washington, DC: U.S. Patent and Trademark Office (2016).
51. Lawrence, J. G. Selfish operons and speciation by gene transfer. *Trends Microbiol.* **5**, 355–359 (1997).
52. Fani, R., Brilli, M. & Lio, P. The origin and evolution of operons: the piecemeal building of the proteobacterial histidine operon. *J. Mol. Evol.* **60**, 378–390 (2005).
53. Lenz, O. & Friedrich, B. A novel multicomponent regulatory system mediates H<sub>2</sub> sensing in *Alcaligenes eutrophus*. *P. Natl. Acad. Sci. USA* **95**, 12474–12479 (1998).
54. Kovach, M. E. *et al.* Four new derivatives of the broad-host-range cloning vector pBBR1MCS, carrying different antibiotic-resistance cassettes. *Gene* **166**, 175–176 (1995).
55. Huang, H. *et al.* Medium optimization and properties of biopolymer Ss from *Sphingomonas sanxanigenens*. *Afr. J. Microbiol. Res.* **6**, 1423–1429 (2012).
56. Delcher, A. L., Harmon, D., Kasif, S., White, O. & Salzberg, S. L. Improved microbial gene identification with GLIMMER. *Nucleic Acids Res.* **27**, 4636–4641 (1999).
57. Lowe, T. M. & Eddy, S. R. tRNAscan-SE: a program for improved detection of transfer RNA genes in genomic sequence. *Nucleic Acids Res.* **25**, 955–964 (1997).
58. Tamura, K., Stecher, G., Peterson, D., Filipowski, A. & Kumar, S. MEGA6: molecular evolutionary genetics analysis version 6.0. *Mol. Biol. Evol.* **30**, 2725–2729 (2013).
59. Langille, M. & Brinkman, F. S. IslandViewer: an integrated interface for computational identification and visualization of genomic islands. *Bioinformatics* **25**, 664–665 (2009).
60. Siguier, P., Perochon, J., Lestrade, L., Mahillon, J. & Chandler, M. ISfinder: the reference centre for bacterial insertion sequences. *Nucleic Acids Res.* **34**, D32–D36 (2006).
61. Levin, J. Z. *et al.* Comprehensive comparative analysis of strand-specific RNA sequencing methods. *Nat. methods* **7**, 709–715 (2010).

62. Trapnell, C. *et al.* Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat. Biotechnol.* **28**, 511–515 (2010).
63. Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinformatics* **12**, 323. (2011).
64. Vijayan, V., Jain, I. H. & O'Shea, E. K. A high resolution map of a cyanobacterial transcriptome. *Genome Biol* **12**, R47 (2011).
65. Noda, S. *et al.* Molecular structures of gellan gum imaged with atomic force microscopy in relation to the rheological behavior in aqueous systems. 1. Gellan gum with various acyl contents in the presence and absence of potassium. *Food Hydrocolloid.* **22**, 1148–1159 (2008).

## Acknowledgements

This work was supported by the National Natural Science Foundation of China (Grant No. 31571790, 41373074) and China Postdoctoral Science Foundation (Grant No. 2016M601251). We also thank Jie Chen (Shanghai Majorbio Bio-pharm Biotechnology Co., Ltd.) for the genome data analysis.

## Author Contributions

M.W. and H.H. contributed equally to this work and wrote the manuscript. T.M., H.H. and G.L. designed all the research. M.W. and H.H. performed the main experiments. Y.R., M.W. and H.H. analyzed the data. Z.S., X.L., X.D., G.G. and M.R. carried out some experiments. All authors reviewed and approved the manuscript.

## Additional Information

**Supplementary information** accompanies this paper at <http://www.nature.com/srep>

**Competing Interests:** The authors declare no competing financial interests.

**How to cite this article:** Wu, M. *et al.* The evolutionary life cycle of the polysaccharide biosynthetic gene cluster based on the *Sphingomonadaceae*. *Sci. Rep.* **7**, 46484; doi: 10.1038/srep46484 (2017).

**Publisher's note:** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



This work is licensed under a Creative Commons Attribution 4.0 International License. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in the credit line; if the material is not included under the Creative Commons license, users will need to obtain permission from the license holder to reproduce the material. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>

© The Author(s) 2017